LaRS: Latent Reasoning Skills for Chain-of-Thought Reasoning

Anonymous ACL submission

Abstract

001

011

012

014

015

021

027

034

037

038

041

Chain-of-thought (CoT) prompting has emerged as a popular in-context learning (ICL) approach for large language models (LLMs), 004 especially when tackling complex reasoning tasks. Traditional ICL approaches construct 006 prompts using examples that contain questions similar to the input question. However, CoT prompting, which includes crucial intermediate reasoning steps (rationales) within its examples, necessitates selecting examples based on these rationales rather than the questions themselves. Existing methods require human experts or pre-trained LLMs to describe the skill, a high-level abstraction of rationales, to guide the selection. These methods, however, are often costly and difficult to scale. Instead, this paper introduces a new approach named Latent Reasoning Skills (LaRS) that employs unsupervised learning to create a latent space representation of rationales, with a latent variable called a *reasoning skill*. Concurrently, LaRS learns a reasoning policy to determine the required reasoning skill for a given question. Then the ICL examples are selected by aligning the reasoning skills between past examples and the question. Our approach is theoretically grounded and sample-efficient, eliminating the need for helper LLM inference or manual prompt design. Empirically, LaRS achieves performance comparable to SOTA rationale-based selection methods, saving thousands of LLM inferences and significantly reducing the time required to process the example bank.

1 Introduction

Large Language Models (LLMs) exhibit remarkable capabilities in solving various downstream tasks through in-context learning (ICL) (Brown et al., 2020), even without being explicitly trained on the distribution of in-context examples (Vaswani et al., 2017; Devlin et al., 2019; Rae et al., 2021;



Figure 1: CoT prompting with examples selected by (a) similar questions and (b) similar skills that (mis)match the skills in their rationales.

Chowdhery et al., 2022; Wei et al., 2022a). Using in-context learning, LLMs generate output for an input query by conditioning on a prompt that contains a few input-output demonstrations.

Reasoning tasks have proven to be particularly difficult for language models and NLP in general (Rae et al., 2021; Bommasani et al., 2021; Nye et al., 2021). In the recent literature, chainof-thought (CoT) prompting, an ICL method, has been proposed to improve LLMs on a wide spectrum of reasoning tasks by guiding LLMs to produce a sequence of intermediate steps (rationale) for generating a (better) final answer (Cobbe et al., 2021a; Wei et al., 2022b; Suzgun et al., 2022). The prompts for CoT are composed of *demonstrations* that contain not only input and output, but also the



Figure 2: An overview of LaRS including a pre-processing stage (left) and a selection stage (right).

rationales for why the output holds.

The core challenge for ICL lies in designing effective demonstrations to prompt LLMs. Much evidence has indicated the significant impact of demonstrations on the performance of ICL (Lu et al., 2021; Liu et al., 2021). To form a prompt, one important setting considers selecting demonstrations from an existing example bank, termed demonstration selection (Dong et al., 2022). While a variety of methods exist in the ICL literature for automating this process, CoT prompts are distinct in that they include not only questions and answers but also specially-designed rationales. This distinction highlights the importance of rationales in selecting demonstrations for CoT prompting. Specifically, CoT prompting should select demonstrations that illustrate relevant skills within their rationales to effectively address a given question. For instance, in solving math word problems (as depicted in Fig. 1), a useful rationale involves computing addition to get the correct answer. Selecting few-shot examples based on the question similarity (Fig. 1a) might lead to examples showcasing subtraction and generate incorrect rationales. However, skill-based selection (Fig. 1b) can align the skills between examples and the given question, which leads to correct answers guided by relevant rationales.

To achieve such a skill-based demonstration selection, An et al. (2023b) introduces Skill-KNN, which employs pre-trained LLMs to generate skill descriptions. Then, the few-shot examples are selected based on the embedding of the skill descriptions computed by another pre-trained embedding model. This process is illustrated by Fig. 2 (left). Although this approach is straightforward, the LLM-generated skill descriptions can be somewhat arbitrary, heavily relying on the manually crafted prompts. This reliance constrains its wider applicability across diverse reasoning tasks. Moreover, the approach requires to generate a unique skill description for each example, which limits its scalability to larger example banks.

098

100

101

102

103

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

Rather than relying on LLMs, we introduce Latent Reasoning Skill Discovery (LaRS), a new skill-based demonstration selection method. This approach learns skills as latent space representations of rationales through unsupervised learning. The essence of LaRS lies in a unique formulation for the generation of rationales, which we term the latent *skill model*. This model, inspired by the principles of topic models (Xie et al., 2021a), conditions the generation of a rationale on both a given question and a latent variable, called a *reasoning skill*. This latent variable embodies a high-level abstraction of the rationales, such as formats, equations, or knowledge.

Under the skill model formulation, LaRS utilizes a Conditional Variational Auto-encoder (CVAE) to approximate the generation of rationales on a small dataset from the example bank. As a result, two probabilistic models can be learned concurrently: (1) a *reasoning skill encoder* that maps an example to the actual reasoning skills demonstrated in the rationale; and (2) a *reasoning policy* that predicts the reasoning skills required for a particular question. This method of learning through a CVAE, especially when applied to a small dataset from the example bank, is both cost-efficient and fast compared to Skill-KNN. Fig. 2 presents an overview of LaRS, including a comparison of its computational efficiency relative to Skill-KNN.

The efficacy of LaRS is evaluated on four different benchmarks based on four backbone LLMs with varying scales. The method is also compared with baseline approaches, including an oracle method that assumes access to ground truth rationales. LaRS achieves comparable performance to

Skill-KNN with no extra LLM inferences and also
matches the oracle performance in almost half of
the experiments. A summary of this paper's contribution is as follows:

- We propose LaRS, a novel unsupervised demonstration selection approach for CoT prompting, and empirically verify its effectiveness through four sets of experiments
- We introduce the latent skill model, a plausible formulation for CoT reasoning, which has illuminated a deeper understanding of CoT prompting.
- We present theoretical analyses of the optimality of the latent-skill-based selection method.

2 Related Work

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

165

166

168

170

171

172

173

174

175

177

178

179

181

This section discusses related work in two different directions.

2.1 CoT Reasoning

CoT prompting is a special prompt design technique that encourages LLMs to generate intermediate rationales that guide them towards providing accurate final answers. These rationales can exhibit remarkable flexibility in their styles. For instance, the original work by (Wei et al., 2022b) specially designs rationales in the in-context demonstrations to suit different reasoning tasks. Moreover, novel prompt designs that highlight diverse formats of the rationales have emerged to enhance CoT prompting. For example, (Kojima et al., 2022) proposed Program of Thoughts (PoT) that disentangles textual reasoning from computation, with the latter specially handled through program generation.

In contrast to manual design, our method LaRS can be thought of as automatic discovery of diverse rationale styles from an example bank. This method can also dynamically select reasoning skills based on the specific questions. Worth noting, (Chen et al., 2023) introduces SKills-in-Context (SKiC), which confines rationale generation to predefined "skills" within the prompt. Although sharing a similar motivation to LaRS, we emphasize two crucial distinctions: (1) while SKiC relies on manual "skills" design, LaRS automatically discovers them, (2) SKiC presents a full list of "skills" in the prompt, allowing LLMs to select from them, whereas LaRS learns the skill selection from the example bank, explicitly instructing LLMs on which skill to employ through in-context examples.

2.2 Demonstration Selection

Demonstration selection refers to a special setting, where the prompts are constructed by selecting examples from an example bank. In this context, our LaRS aligns with the paradigm of unsupervised demonstration selection, which involves designing heuristics for this selection process. A variety of heuristics have been explored, including similarity (Gao et al., 2021; Hu et al., 2022), diversity (Zhang et al., 2022), coverage (Gupta et al., 2023), and uncertainty (Diao et al., 2023). Among these, Skill-KNN ((An et al., 2023b)) shares the closest resemblance to our approach. However, Skill-KNN relies on pre-trained LLMs to provide "skill" annotations, which could be arbitrary and resource-intensive, requiring extensive inferences of LLMs and human prompt design. In contrast, LaRS automatically discovers reasoning skills by learning a very lightweight CVAE. In addition, the selections based on these discovered reasoning skills are theoretically-grounded based on the latent skill model and the theoretical analyses presented in this paper.

182

183

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

226

3 Formulation

In this section, we formally describe the *skill model*, a new formulation for explaining the generation of rationales in CoT reasoning. In Section 3.1, the skill model is first introduced to describe the human-generated rationales. Then, Section 3.2 illustrates how the skill model can be adapted to LLM-generated rationales. Finally, leveraging the concept of reasoning skill as outlined in the skill model, a new latent-skill-based demonstration selection method is formally described in Section 3.3.

3.1 Skill Model

Let \mathcal{X} be the set of all sequences of tokens, \mathcal{Z} be the continuous vector space of latent reasoning skills, and P_H denotes the probability distribution of real-world natural language. CoT reasoning is to generate a rationale $R \in \mathcal{X}$ given a question $Q \in \mathcal{X}$, whose correctness¹ can be verified by an indicator function $\mathbb{1}(R, Q) :=$ $\mathbb{1}(R$ is the correct rationale for Q).

The skill model assumes that the real-world conditional distribution of R given Q can be described

¹For math word problems, whose answers are discrete labels, the correct rationale should contain the correct answer label as the final step. For code generation, the correct rationale should be the correct code.

as follows:

227

233

235

238

239

240

241

242

243

244

245

247

254

256

257

260

261

264

265

267

271

272

274

$$P_H(R \mid Q) = \int_{\mathcal{Z}} P_H(R \mid z, Q) P_H(z \mid Q) dz \quad (1)$$

where, $P_H(z \mid Q)$ is the posterior of selecting latent reasoning skills in human reasoning, called a reasoning policy. $P_H(R \mid z, Q)$ is the posterior distribution of generating R given a question Q and a reasoning skill z. A causal graph illustrating such a generation process involving a latent reasoning skill z is presented in Fig. 3 on the left.

Unlike (Wang et al., 2023), this formulation considers a dependency of z on Q reflecting a preference for selecting particular reasoning skills to solve a given question. We justify this formulation as follows:

- 1. Rationales can exhibit remarkable flexibility, manifesting diverse formats, topics, and knowledge, which can naturally be abstracted into the high-level concepts of reasoning skills.
- 2. The selection of these skills is not bound by strict determinism. For instance, diverse reasoning paths and formats could all contribute toward finding the correct final answer. Therefore, real-world data is a mixture of diverse skills captured by a stochastic reasoning policy $P_H(z \mid Q)$.

3.2 CoT prompting

LLMs are pre-trained conditional generators. Given an input query $X \in \mathcal{X}$, the conditional distribution of an output $Y \in \mathcal{X}$ generated by LLMs can be written as $P_M(Y \mid X)$. LLMs are usually trained on generic real-world data distribution such that $P_M(Y \mid X) \approx P_H(Y \mid X)$.

Prior studies have presented an implicit topic model formulation in explaining the in-context learning mechanisms of LLMs (Wang et al., 2023; Xie et al., 2021a). Similarly, we posit that LLMs can be viewed as implicit skill models for generating rationales. To elaborate, when generating rationales, LLMs' conditional distribution $P_M(R \mid Q)$ can be extended as follows (with illustrations in Fig. 3 on the left):

$$P_M(R \mid Q) = \int_{\mathcal{Z}} P_M(R \mid z, Q) P_M(z \mid Q) dz \quad (2)$$

This implicit skill model assumes that LLMs also infer reasoning skills *z*, which resembles the realworld generation of rationales.

The above formulation only encompasses the zero-shot generation of rationales. In practice,



Figure 3: Causal graphs for prompting with zero-shot/human (left), zero-shot CoT (middle), and few-shot CoT (right) for generating rationales via skills. The dashed arrow from Q to z indicates possible sub-optimal inference of the reasoning skills from both human and zero-shot LLM generations.

prompts are commonly provided to guide LLMs' generation. In general, two CoT prompting strategies exist: zero-shot CoT, employing a prompt comprising a short prefix and a test question, and fewshot CoT, employing a prompt containing pairs of questions and rationales. Denoting $pt \in \mathcal{X}$ as a prompt, a unified formulation for both prompting strategies can be derived as follows:

$$P_M(R \mid pt) = \int_{\mathcal{Z}} P_M(R \mid z, Q) P_M(z \mid pt) dz \quad (3)$$

0-shot CoT: $pt = (\text{prefix}, Q) \text{ or } (Q, \text{prefix})$
 k -shot CoT: $pt = (Q_1, R_1, \cdots, Q_k, R_k, Q)$

Here, the formulation is simplified such that the use of prompts only influences the probability distribution of z. For instance, a prefix specifying the generation's format can be interpreted as specifying the reasoning skill z by shaping the distribution from $P_M(z \mid Q)$ to $P_M(z \mid pt)$. This simplification aligns with empirical evidence suggesting that in-context examples serve as mere pointers to retrieve already-learned knowledge within LLMs (Shin et al., 2020; Min et al., 2022).

Drawing upon this formulation, we can gain insight into the failure of zero-shot generation. In general, real-world data is inherently noisy, indicating that the reasoning policy $P_H(z \mid Q)$ may be sub-optimal, and the reasoning skills are not chosen to maximize the accuracy of answering a test question. Trained on this generic real-world data distribution, $P_M(z \mid Q)$ could also be sub-optimal, leading to the failure of zero-shot generation. On the other hand, CoT prompting improves the reasoning performance by shaping the distribution of reasoning skills using carefully-designed prompts that contain either instructions or few-shot examples.

3.3 Skill-Based Demonstration Selection

The analysis above suggests that the key to the success of CoT prompting is to design an effective prompt that shapes the posterior distribution of 306

307

308

310

311

312

313

381

382

383

384

385

387

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

361

362

335

331

332

333

315

316

317

320

321

324

325

337

339 340

341

343

352

360

348

349

347

345

346

4

examples that align with $P_E(z \mid Q)$ are intuitively

Method

illustrated in Figure 2.

is given by

 $(R,Q) \in \mathcal{D}_E.$

according to $P_E(z \mid Q)$. Since the example bank is usually specially-crafted and contains rationales showcasing "better" reasoning skills, the in-context

tions of the example bank and LLMs.

more effective. In Section 4.3, we provide theoreti-

cal analysis of the optimality of this skill-based se-

lection when conditioned on certain ideal assump-

To enable the skill-based demonstration selec-

tion (Definition 1), we introduce our approach

LaRS, which involves learning a conditional vari-

ational autoencoder (CVAE) to approximate P_E

using the data from the example bank \mathcal{D}_E . We

then outline a practical demonstration selection

process aligning with the skill-based selection. The

schematic overview of LaRS (right) and the corre-

sponding demonstration selection process (left) are

reasoning skills, assuming that the real-world dis-

tribution $P_H(z \mid Q)$ is potentially sub-optimal. In

contrast to the real-world distribution, the demon-

stration selection problem assumes access to an

example bank of question-rationale pairs, denoted

as $\mathcal{D}_E = \{(R, Q)\}$. This example bank is usu-

ally specially-crafted and has a distribution dif-

ferent from the real-world distribution. Denot-

ing P_E as the distribution of the example bank,

R is distributed according to $P_E(R \mid Q)$ for all

Given \mathcal{D}_E , the demonstration selection is to se-

lect a few question-rationale pairs from \mathcal{D}_E . As-

suming that each selected demonstration is i.i.d, a

demonstration selection method can be uniquely

defined as a probabilistic model $g(Q, R|Q_{\text{test}}) :=$

 $\mathcal{X} \mapsto \Delta(\mathcal{X})$ that maps a test question Q_{test} to a

probability distribution of demonstrations. Then,

we can formally define the skill-based demonstra-

Definition 1 Skill-based demonstration selection

 $g_{skill}(Q, R \mid Q_{test}) = \int_{\mathcal{Z}} P_E(Q, R \mid z) P_E(z \mid Q_{test}) dz$

tion selection method as follows:

Intuitively, this selection method maximizes the probability of a selected demonstration showcasing the reasoning skill that is likely to be chosen 4.1 Latent Reasoning Skill Discovery

The conditional variational autoencoder (CVAE) has emerged as a popular approach for modeling probabilistic conditional generation. As one specific case, the skill model, introduced in this paper, can effectively be represented as a CVAE. Therefore, we introduce LaRS that employs a CVAE to approximate the generation of rationales using the data from the example bank $\mathcal{D}_E = \{(Q, R)\}.$

In particular, this CVAE includes three coupled models: an encoder model, a decoder model, and a reasoning policy model, independently parameterized by ω , ψ , and ϕ respectively. Drawing from the notations introduced in the skill model, the reasoning policy model is a conditional Bayesian network $\pi_{\phi}(z \mid Q)$, determining the posterior distribution of latent reasoning skill z given a question Q. The decoder model is also a conditional Bayesian network $p_{\psi}(R \mid z, Q)$ that generates a rationale R, conditioned on both Q and z, where z is sampled from $\pi_{\phi}(z \mid Q)$. Finally, the encoder model $q_{\omega}(z \mid Q, R)$ is another conditional Bayesian network, mapping a question-rationale pair to z. In this paper, we train this CVAE using classical variational expectation maximization and the reparameterization trick.

Specifically, the classical variational expectation maximization optimizes a loss function as follows:

 $\mathcal{L}_{\text{CVAE}}(\phi, \omega, \psi) = \mathcal{L}_{\text{recon}} + \mathcal{L}_{\text{KL}}$ (4) $\mathcal{L}_{\text{recon}} = -\mathbb{E}_{(Q,R)\sim\mathcal{D}_E, z\sim q_\omega(|Q,R)}[\log p_{\psi}(R|z,Q)]$ $\mathcal{L}_{\mathrm{KL}} = \mathbb{E}_{(Q,R)\sim\mathcal{D}_E}[\mathbf{D}_{\mathrm{KL}}(q_\omega(z \mid Q, R) \parallel \pi_\phi(z \mid Q))]$

By training to minimize this loss function, q_{ω} and π_{ϕ} can be learned to effectively approximate the conditional distributions $P_E(z \mid Q, R)$ and $P_E(z \mid Q)$. It is worth noting that the decoder model acts an auxiliary model that only roughly reconstructs rationales for the purpose of training the encoder and the reasoning policy model, and is not deployed to generate rationales in the downstream tasks.

Ideally, all three models would be represented by language models, processing token sequences as input and generating token sequences as output. However, training full language models for demonstration selections can be computationally expensive. Instead, we adopt a pre-trained embedding model denoted as $f : \mathcal{X} \mapsto \Theta$, which maps the token space \mathcal{X} to an embedding

Algorithm 1 Demonstration selection

Input: Test question Q_{test} , a pre-trained embedding model f, a reasoning policy $\pi_{\phi}(z|f(Q))$, a reasoning skill encoder $q_{\omega}(z|f(Q,R))$, and an example bank $\mathcal{D}_E = \{(Q^j, R^j)\}_j$.

Parameter: shot number k

Output: $(Q_1, R_1, Q_2, R_2, \cdots, Q_k, R_k)$

- 1: Compute $z_{\text{test}} \leftarrow \text{mean of } \pi(z|f(Q_{\text{test}}))$
- 2: for each (Q^j, R^j) in \mathcal{D}_E do
- 3: Compute $z_{\text{post}}^j \leftarrow \text{mean of } q_\omega(z|f(Q^j, R^j))$

4: Compute
$$r^j = \frac{z_{\text{test}} \cdot z_{\text{post}}^{j-1}}{|z_{\text{test}}| \cdot |z_{\text{post}}^j|}$$

- 5: end for
- 6: Select top-k demonstrations with the largest r^j and sort them in ascending order, denoted as (Q₁, R₁, Q₂, R₂, ..., Q_k, R_k).
 7. return (Q, B, Q, B), Q
- 7: return $(Q_1, R_1, Q_2, R_2, \cdots, Q_k, R_k) = 0$

space Θ . Consequently, the decoder model, encoder model, and reasoning policy model transform into $p_{\psi}(f(R)|z, f(Q)), q_{\omega}(z|f(Q, R))$, and $\pi_{\phi}(z|f(Q))$, respectively. They now condition on and generate the embeddings instead of the original tokens. In the actual implementation, we use the same feed-forward neural network to represent both π_{ϕ} and q_{ω} , predicting the mean and variance of Gaussian distributions of latent reasoning skills. On the other hand, p_{ψ} is a feed-forward neural network that deterministically predicts a value in the embedding space.

4.2 Demonstration Selection

Since the distribution $P_E(Q, R | z)$ in Definition 1 is practically intractable, we propose a selection process that effectively aligns with the skill-based selection using the learned π_{ϕ} and q_{ω} . For a given test question Q_{test} , the desirable reasoning skill $z_{\text{test}} = \arg \max_z [\pi_{\phi}(z|f(Q_{\text{test}}))]$ can be computed using the reasoning policy. Subsequently, each example from the example bank can be scored based on the cosine similarity between z_{test} and z_{post} , where $z_{\text{post}} = \arg \max_z [q_{\omega}(z|Q, R))]$ represents the maximum likelihood skill of the current example. Finally, a CoT prompt can be constructed by selecting the top-k examples according to the computed scores. The step-by-step procedure is outlined in Algorithm 1.

4.3 Theoretical Analysis

In this section, we provide a theoretical analysis of the optimality of the skill-based selection by Definition 1.

Let $P_M(R \mid Q, g)$ denotes LLMs' conditional distribution of a rationale R given a test question Q under a demonstration selection method $g. P_M(R \mid Q, g)$ can be extended as follows: 442

$$P_M(R \mid Q, g) = \int_{\mathcal{X}^k} P_M(R \mid pt) \prod_{i=1}^k [g(Q_i, R_i \mid Q)d(Q_i, R_i)]$$

$$44$$

Here, each demonstrations (Q_i, R_i) is independently sampled from $g(Q_i, R_i | Q), \forall i = 1, \dots, k$. These k demonstrations form a prompt $pt = (Q_1, R_1, \dots, Q_k, R_k, Q)$.

We want to show that $P_M(R \mid Q, g)$ is the optimal conditional distribution that maximizes the accuracy of rationales if the selection follows skillbased selection method or $g = g_{skill}$. We begin by defining the optimal conditional distribution as follows:

Definition 2 Optimal conditional distribution of rationales given questions $P^*(R \mid Q)$ is given by:

$$P^*(R \mid Q) = \underset{P(\cdot \mid Q) \in \Delta(\mathcal{X})}{\operatorname{arg\,max}} \int_{\mathcal{X}} \mathbb{1}(R, Q) P(R \mid Q) dR$$

Here $\mathbb{1}(R,Q)$ is the indicator function of the correctness of R given a question Q (see Section 3.1).

Then, we state two major assumptions as follows:

Assumption 1 *Example bank is sampled from the* optimal conditional distribution, or $P_E(R \mid Q) = P^*(R \mid Q)$.

Assumption 2 Humans and LLMs are expert rationale generators given reasoning skills and questions, meaning that

$$P_H(R \mid z, Q) = P_E(R \mid z, Q) = P_M(R \mid z, Q).$$

Assumption 1 is rooted in the fact that example banks are human-crafted that contains the most useful rationales for answering the questions. In Assumption 2, P_M capturing P_H is a common assumption in the literature studying LLMs (Xie et al., 2021b; Saunshi et al., 2020; Wei et al., 2021). $P_E(R \mid z, Q) = P_H(R \mid z, Q)$ is based on the assumption that reasoning skills are shared across humans, and the generation of rationales is identical given the same reasoning skills and questions.

Based on the above definiton and two assumptions, we prove the following theorem.

428

429

430

431

432

433

434

435

436

437

407

408

409

438

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

566

567

568

521

Theorem 1 A LLM gives the optimal conditional distribution of rationales given questions:

$$P_M(R \mid Q, g_{skill}) = P^*(R \mid Q)$$

If (1) it is prompted by $k \to \infty$ in-context examples selected by the skill-based selection g_{skill} defined by Definition 1, (2) Assumption 2 and Assumption 1 hold.

Appendix D presents the proof for Theorem 1.

5 Experiments

480

481

482

483

484

485

486

487

488

489

491

492

493

494

495

496

497

498

499

502

504

508

510

512

513

514

515

516

518

519

520

This section describes the experiment settings, including benchmarks, compared selection methods, backbone models, and hyper-parameters. Lastly, the main results of these experiments are presented.

5.1 Dataset

For benchmarking, the selection methods are evaluated on four challenging datasets, including two datasets of Math Word Problem (MWP): **TabMWP**, **GSM8K**, one text-to-SQL dataset: **Spider**, and one semantic parsing dataset: **COGS**.

Each dataset is split into a training set used to learn LaRS models and a test set used to evaluate the selection methods. While the training sets may potentially be large, we use randomly sampled 1K examples from the training set as the example bank, from which, the examples can be selected for CoT prompting. Detailed descriptions of the datasets and splitting are presented in Appendix A.

To measure the performances, we use the answer accuracy for **TabMWP** and **GSM8K**, with the answers extracted by searching the texts right after a prefix The answer is. For **Spider**, we use the official execution-with-values accuracy². For **COGS**, we report the exact-match accuracy for semantic parsing.

5.2 Selection Methods

Our method LaRS is compared with the following four baselines. All the hyper-parameters related to these methods are listed in Appendix A.

Skill-KNN This baseline is a SOTA skill-based selection method that uses a pre-trained LLM (gpt-3.5-turbo) to generate skill descriptions for both questions in the example bank and the test question. Then, the method selects examples with the most similar skill descriptions as the test question skill description to form the prompt. The cosine similarity is computed based a pre-trained embedding model.

Random This baseline randomly selects k incontext examples from the example bank. For each test question, the accuracy is reported as an average over three independent random selections.

Retrieval-Q This baseline employs a pre-trained embedding model to encode a test question, and selects in-context examples based on the cosine similarity between embeddings from examples' questions and the test question.

Retrieval-R (oracle) This baseline employs a pre-trained embedding model to encode the ground-truth rationale of a test question, and selects incontext examples based on the cosine similarity between examples' rationales and the ground-truth rationale.

5.3 Backbones and Hyper-parameters

In terms of the backbone models, the ICL is conducted by two OpenAI language models: text-davinci-003 and gpt-3.5-turbo, one Anthropic model: Claude-v2, and one smaller-scale Falcon-40B-Instruct (Xu et al., 2023). All the embedding is computed by a pre-trained embedding model, Deberta-v2-xlarge (He et al., 2021). We also investigate different choices of embedding model in Section B.

During inference, the temperature is set to 0 (i.e., greedy decoding) to reduce the variance. The CoT prompts contain k = 2, 4, 4, 8 in-context examples for **TabMWP**, **GSM8K**, **Spider**, and **COGS**, respectively.

5.4 Performance comparison results

Table 1 presents a summary of the experimental results. Detailed descriptions are provided as follows:

LaRS is comparable to Skill-KNN. Across all four benchmarks and three backbone models tested, our proposed LaRS outperforms Skill-KNN in 7 out of 12 experiments. This result underlines the effectiveness of the latent reasoning skills learned through unsupervised learning with small CVAE models, achieving comparable performance to the skill descriptions crafted by extensive pre-trained LLMs. Furthermore, LaRS consistently achieved superior results compared to Retrieval-Q, which uti-

²We use the official evaluation scripts for Spider in https://github.com/taoyds/test-suite-sql-eval.

Method	TabMWP	GSM8K	Spider	COGS	
Backbone: gpt-3.5-turbo					
Random	62.4 _{+0.0}	75.7 +0.0	46.8 +0.0	67.5 _{+0.0}	
Retrieval-Q	72.3 +9.9	75.6 _{-0.1}	49.9 _{+3.1}	88.5 _{+21.0}	
Skill-KNN	78.3 +15.9	$75.0_{-0.7}$	58.4 +11.6	94.6 +27.2	
LaRS (ours)	$78.1_{+15.7}$	76.8 +1.1	53.0 _{+6.2}	94.6 +27.1	
Retrieval-R (oracle)	77.4 _{+15.0}	$75.5_{-0.2}$	$64.4_{+17.6}$	95.7 _{+28.2}	
Backbone: text-davinci-003					
Random	69.3 _{+0.0}	62.2 +0.0	47.1 +0.0	73.4 +0.0	
Retrieval-Q	76.5 +7.2	$62.7_{+0.5}$	50.2 +2.9	92.1 _{+18.7}	
Skill-KNN	80.6 +11.3	62.0 _{-0.2}	56.3 +9.8	96.8 _{+23.4}	
LaRS (ours)	80.8 +11.5	62.7 +0.5	$48.6_{+1.5}$	96.6 +23.2	
Retrieval-R (oracle)	80.4 +11.1	$63.8_{\ +1.6}$	$67.3_{+20.2}$	97.3 _{+23.9}	
Backbone: Claude-v2					
Random	77.7 _{+0.0}	86.9 +0.0	40.2 +0.0	77.6 _{+0.0}	
Retrieval-Q	$80.1_{+2.4}$	88.2 +1.3	45.5 +5.3	93.5 _{+15.9}	
LaRS (ours)	80.9 +3.2	88.3 +1.4	47.7 +7.5	96.6 +19.0	
Retrieval-R (oracle)	80.3 +2.6	$88.4_{+1.5}$	60.8 _{+20.6}	$97.3_{\ +19.7}$	
Backbone: Falcon-40B-Instruct					
Random	45.7 +0.0	38.8 +0.0	20.6 +0.0	45.1 +0.0	
Retrieval-Q	51.9 _{+6.2}	37.3 _{-1.5}	$22.1_{+1.5}$	$73.9_{+28.8}$	
Skill-KNN	55.9 _{+10.2}	40.3 +1.5	23.7 +2.9	81.0 +35.9	
LaRS (ours)	57.7 +12.0	$39.1_{+0.3}$	24.8 +4.2	89.5 +44.4	
Retrieval-R (oracle)	$61.2_{+15.5}$	$40.4_{+1.6}$	39.9 +19.3	90.3 +45.2	

Table 1: Main results (%) across all backbone models and datasets. Numbers in **bold** represent the best results for each backbone model across all selection methods. The subscripted gray values indicate the relative improvement over Random selection.

lizes raw question embeddings. This observation underscores that the success of LaRS is attributed to the learned reasoning skill representation rather than solely relying on the raw information from the questions. As depicted in Figure 1 within Appendix B, we manually identified 12 skill labels from **TabMWP**. The scatter plots illustrate the distinct separation of these skills by LaRS, Skill-KNN, question embedding, and rationale embedding, respectively.

LaRS is LM-agnostic. The superiority of LaRS is consistent across four different LMs, including the both open-source and proprietary models, despite not being specifically trained for any of these LMs. This finding underscores the universality of the learned latent reasoning skill representation, enabling any LMs to benefit from it.

586LaRS is computationally efficient. In Table 2,587we present a comparison of computational over-588head, including computing time, estimated cost589for pre-processing the example bank, and cost for590each input query during testing, among Retrieval-Q,591LaRS, Skill-KNN, and the supervised demonstra-592tion selection method PromptPG (Lu et al., 2022).593These estimates are based on the TabMWP dataset,

	Accuracy (%)	Time (h)	Cost (training)	Cost per query
Retrieval-Q	72.3	0	\$0	\$0.02 +%0
LaRS (ours)	78.1	0.5	\$0	\$0.02 +%0
Skill-KNN	78.3	2	\$30	\$0.05 _{+%150}
PromptPG	74.2	6	\$50	$0.02_{+\%0}$

Table 2: Comparison of accuracy and computational overhead, including computing time, estimated cost for pre-processing the example bank, and average cost per input query, among four selection methods on the **TabMWP** dataset with an example bank of size 1000. The grey percentages represent the increased cost ratio associated with each selection method.

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

using GPT-4 as the backbone model, with the API pricing at \$0.03 per 1k tokens. Retrieval-Q incurs the lowest computational overhead but exhibits the poorest performance among the methods evaluated. Our method achieves accuracy comparable to Skill-KNN, requiring 1000 fewer LLM inferences (approximately \$30 savings) and reducing computing time by 1.5 hours during training, along with more than 100% less cost per input query. Remarkably, even though it uses more computational resources, PromptPG does not outperform our unsupervised method in terms of accuracy.

6 Conclusions

This paper introduces LaRS, a novel demonstration selection method designed for CoT prompting. LaRS bases the selection on reasoning skills, which are latent representations discovered by unsupervised learning from rationales via a CVAE. Based on the experiments conducted across four LLMs and over four different reasoning tasks, LaRS manifests comparable performance on selecting effective few-shot examples for CoT reasoning while requiring no extra LLM inference and saving hours in pre-processing the example bank.

7 Limitations

Despite the success of LaRS, a few limitations and potential future directions are worth noting. First, the impact of the order of examples in the prompts is not considered. Introducing additional heuristics to sort the examples could potentially lead to better performances. Second, in the CVAE, the decoder is represented by an MLP neural network. However, it would be ideal to represent the decoder as a prompt-tuning module, which aligns better with the implicit skill model assumption. Finally, one single reasoning skill might not be sufficient to represent the entire rationale that might contain multiple steps of reasoning. Learning and selecting reasoning skills for each individual reasoning step is an interesting direction to explore.

585

569

References

634

635

639

641

643

653

664

668

675

676

677

679

690

- Shengnan An, Zeqi Lin, Qiang Fu, B. Chen, Nanning Zheng, Jian-Guang Lou, and D. Zhang. 2023a. How do in-context examples affect compositional generalization? *ArXiv*, abs/2305.04835.
 - Shengnan An, Bo Zhou, Zeqi Lin, Qiang Fu, Bei Chen, Nanning Zheng, Weizhu Chen, and Jian-Guang Lou.
 2023b. Skill-based few-shot selection for in-context learning. *arXiv preprint arXiv:2305.14210*.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, S. Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen A. Creel, Jared Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren E. Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas F. Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, O. Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Benjamin Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, J. F. Nyarko, Giray Ogut, Laurel J. Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Robert Reich, Hongyu Ren, Frieda Rong, Yusuf H. Roohani, Camilo Ruiz, Jack Ryan, Christopher R'e, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishna Parasuram Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei A. Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2021. On the opportunities and risks of foundation models. ArXiv, abs/2108.07258.
 - Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, T. J. Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeff Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. ArXiv, abs/2005.14165.
 - Jiaao Chen, Xiaoman Pan, Dian Yu, Kaiqiang Song, Xiaoyang Wang, Dong Yu, and Jianshu Chen. 2023.

Skills-in-context prompting: Unlocking compositionality in large language models. *arXiv preprint arXiv:2308.00304*.

695

696

697

698

699

700

701

703

704

705

706

707

709

710

711

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

737

738

739

740

741

742

743

744

745

746

747

749

750

- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam M. Shazeer, Vinodkumar Prabhakaran, Emily Reif, Nan Du, Benton C. Hutchinson, Reiner Pope, James Bradbury, Jacob Austin, Michael Isard, Guy Gur-Ari, Pengcheng Yin, Toju Duke, Anselm Levskaya, Sanjay Ghemawat, Sunipa Dev, Henryk Michalewski, Xavier García, Vedant Misra, Kevin Robinson, Liam Fedus, Denny Zhou, Daphne Ippolito, David Luan, Hyeontaek Lim, Barret Zoph, Alexander Spiridonov, Ryan Sepassi, David Dohan, Shivani Agrawal, Mark Omernick, Andrew M. Dai, Thanumalayan Sankaranarayana Pillai, Marie Pellat, Aitor Lewkowycz, Erica Moreira, Rewon Child, Oleksandr Polozov, Katherine Lee, Zongwei Zhou, Xuezhi Wang, Brennan Saeta, Mark Díaz, Orhan Firat, Michele Catasta, Jason Wei, Kathleen S. Meier-Hellstern, Douglas Eck, Jeff Dean, Slav Petrov, and Noah Fiedel. 2022. Palm: Scaling language modeling with pathways. ArXiv, abs/2204.02311.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. 2021a. Training verifiers to solve math word problems. *ArXiv*, abs/2110.14168.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, et al. 2021b. Training verifiers to solve math word problems. *arXiv preprint arXiv:2110.14168*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. *ArXiv*, abs/1810.04805.
- Shizhe Diao, Pengcheng Wang, Yong Lin, and Tong Zhang. 2023. Active prompting with chainof-thought for large language models. *ArXiv*, abs/2302.12246.
- Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Zhiyong Wu, Baobao Chang, Xu Sun, Jingjing Xu, and Zhifang Sui. 2022. A survey for in-context learning. *arXiv preprint arXiv:2301.00234*.
- Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making pre-trained language models better few-shot learners. *ArXiv*, abs/2012.15723.
- Shivanshu Gupta, Sameer Singh, and Matt Gardner. 2023. Coverage-based example selection for incontext learning. *ArXiv*, abs/2305.14907.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. Deberta: Decoding-enhanced

752 753	bert with disentangled attention. In International Conference on Learning Representations.
754	Yushi Hu, Chia-Hsuan Lee, Tianbao Xie, Tao Yu,
755	Noah A. Smith, and Mari Ostendorf. 2022. In-
756	context learning for few-shot dialogue state tracking.
757	<i>ArXiv</i> , abs/2203.08568.
758	Najoung Kim and Tal Linzen. 2020. Cogs: A compo-
759	sitional generalization challenge based on semantic
760	interpretation. <i>ArXiv</i> , abs/2010.05465.
761	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yu-
762	taka Matsuo, and Yusuke Iwasawa. 2022. Large lan-
763	guage models are zero-shot reasoners. <i>Advances in</i>
764	<i>neural information processing systems</i> , 35:22199–
765	22213.
766	Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan,
767	Lawrence Carin, and Weizhu Chen. 2021. What
768	makes good in-context examples for gpt-3? In Work-
769	shop on Knowledge Extraction and Integration for
770	Deep Learning Architectures; Deep Learning Inside
771	Out.
772	Pan Lu, Liang Qiu, Kai-Wei Chang, Ying Nian Wu,
773	Song-Chun Zhu, Tanmay Rajpurohit, Peter Clark,
774	and A. Kalyan. 2022. Dynamic prompt learning
775	via policy gradient for semi-structured mathematical
776	reasoning. ArXiv, abs/2209.14610.
777	Yao Lu, Max Bartolo, Alastair Moore, Sebastian
778	Riedel, and Pontus Stenetorp. 2021. Fantastically
779	ordered prompts and where to find them: Over-
780	coming few-shot prompt order sensitivity. <i>ArXiv</i> ,
781	abs/2104.08786.
782	Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe,
783	Mike Lewis, Hannaneh Hajishirzi, and Luke Zettle-
784	moyer. 2022. Rethinking the role of demonstra-
785	tions: What makes in-context learning work? <i>arXiv</i>
786	<i>preprint arXiv:2202.12837</i> .
787	Arvind Neelakantan, Tao Xu, Raul Puri, Alec Rad-
788	ford, Jesse Michael Han, Jerry Tworek, Qiming Yuan,
789	Nikolas Tezak, Jong Wook Kim, Chris Hallacy, et al.
790	2022. Text and code embeddings by contrastive pre-
791	training. arXiv preprint arXiv:2201.10005.
792	 Maxwell Nye, Anders Andreassen, Guy Gur-Ari,
793	Henryk Michalewski, Jacob Austin, David Bieber,
794	David Dohan, Aitor Lewkowycz, Maarten Bosma,
795	David Luan, Charles Sutton, and Augustus Odena.
796	2021. Show your work: Scratchpads for interme-
797	diate computation with language models. <i>ArXiv</i> ,
798	abs/2112.00114.
799	Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie
800	Millican, Jordan Hoffmann, Francis Song, John
801	Aslanides, Sarah Henderson, Roman Ring, Susan-
802	nah Young, Eliza Rutherford, Tom Hennigan, Ja-
803	cob Menick, Albin Cassirer, Richard Powell, George
804	van den Driessche, Lisa Anne Hendricks, Mari-
805	beth Rauh, Po-Sen Huang, Amelia Glaese, Johannes
806	Welbl, Sumanth Dathathri, Saffron Huang, Jonathan
807	Uesato, John F. J. Mellor, Irina Higgins, Antonia

Creswell, Nathan McAleese, Amy Wu, Erich Elsen, Siddhant M. Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, L. Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsimpoukelli, N. K. Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Tobias Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d'Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew G. Johnson, Blake A. Hechtman, Laura Weidinger, Iason Gabriel, William S. Isaac, Edward Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem W. Ayoub, Jeff Stanway, L. L. Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. 2021. Scaling language models: Methods, analysis & insights from training gopher. ArXiv, abs/2112.11446.

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

- Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. arXiv preprint arXiv:1908.10084.
- Nikunj Saunshi, Sadhika Malladi, and Sanjeev Arora. 2020. A mathematical exploration of why language models help solve downstream tasks. ArXiv, abs/2010.03648.
- Taylor Shin, Yasaman Razeghi, Robert L Logan IV, Eric Wallace, and Sameer Singh. 2020. Eliciting knowledge from language models using automatically generated prompts. In Conference on Empirical Methods in Natural Language Processing.
- Mirac Suzgun, Nathan Scales, Nathanael Scharli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc V. Le, Ed Huai hsin Chi, Denny Zhou, and Jason Wei. 2022. Challenging big-bench tasks and whether chain-of-thought can solve them. In Annual Meeting of the Association for Computational Linguistics.
- Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In NIPS.
- Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan Sun. 2022. Towards understanding chain-of-thought prompting: An empirical study of what matters. arXiv preprint arXiv:2212.10001.
- Xinyi Wang, Wanrong Zhu, and William Yang Wang. 2023. Large language models are implicitly topic models: Explaining and finding good demonstrations for in-context learning. arXiv preprint arXiv:2301.11916.
- Colin Wei, Sang Michael Xie, and Tengyu Ma. 2021. Why do pretrained language models help in downstream tasks? an analysis of head and prompt tuning. ArXiv, abs/2106.09226.

Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed Huai hsin Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. 2022a. Emergent abilities of large language models. *Trans. Mach. Learn. Res.*, 2022.

865

866

868

873 874

875

876

878

879

881

882

884

887

889

890

895

897

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed Huai hsin Chi, F. Xia, Quoc Le, and Denny Zhou. 2022b. Chain of thought prompting elicits reasoning in large language models. *ArXiv*, abs/2201.11903.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021a. An explanation of in-context learning as implicit bayesian inference. *ArXiv*, abs/2111.02080.
- Sang Michael Xie, Aditi Raghunathan, Percy Liang, and Tengyu Ma. 2021b. An explanation of incontext learning as implicit bayesian inference. *arXiv preprint arXiv:2111.02080*.
- Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. 2023. Baize: An open-source chat model with parameter-efficient tuning on self-chat data. *arXiv preprint arXiv:2304.01196*.
- Tao Yu, Rui Zhang, Kai-Chou Yang, Michihiro Yasunaga, Dongxu Wang, Zifan Li, James Ma, Irene Z Li, Qingning Yao, Shanelle Roman, Zilin Zhang, and Dragomir R. Radev. 2018. Spider: A largescale human-labeled dataset for complex and crossdomain semantic parsing and text-to-sql task. *ArXiv*, abs/1809.08887.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alexander J. Smola. 2022. Automatic chain of thought prompting in large language models. *ArXiv*, abs/2210.03493.

901

902

903

904

906

907

908

917

918

919

922

924

925

926

931

936

938

942

Appendix: LaRS: Latent Reasoning Skill for Chain-of-Thought Reasoning

A Experimental Details

A.1 Dataset

We provide detailed description of the dataset and the split of train and test set as follows:

TabMWP (Lu et al., 2022) This dataset consists of semi-structured mathematical reasoning problems, comprising 38,431 open-domain grade-level problems that require mathematical reasoning on both textual and tabular data. We use the train set, containing 23,059 examples, to train our LaRS models, and test1k set containing 1K examples to evaluate the selection methods.

Spider (Yu et al., 2018) Spider is a large-scale text-to-SQL dataset. It includes a train set with 7,000
examples and a dev set with 1,034 examples. We use the train set to train our LaRS models, and the dev
set as the test set to evaluate the selection methods.

COGS (Kim and Linzen, 2020) is a synthetic benchmark for testing compositional generalization in semantic parsing. We transform the output format in the same way as An et al. (2023a), and consider a mixture of two sub-tasks: primitive substitution (P.S.) and primitive structural alternation (P.A.). This
 results in a train set of 6916 examples to train our LaRS models and a test set of 1000 examples to evaluate the selection method.

GSM8k (Cobbe et al., 2021b) GSM8k is a dataset containing 8.5K high-quality, linguistically diverse grade school math word problems. It includes a train set of 7.5K problems and a test set of 1319 problems. We use the train set to train our LaRS models, and the test set to evaluate the selection methods.

A.2 Hyper-parameters

LaRS contains a encoder, a decoder, and a reasoning policy model. The reasoning skill is represented as a 128-dimensional continuous space. Both the encoder and the reasoning policy model are represented as a feed-forward multiple layer perception (MLP) with two 256-unit hidden layers, predicting the mean and variance of a multivariate Gaussian distribution in the latent space of reasoning skills. The decoder is a MLP with two 256-unit hidden layers that predicts a value in the embedding space deterministically. The dimension of the embedding space depends on the choice of pre-trained embedding models. The models are trained using the loss function in Equation 4 with a batch size of 256 and a learning rate of 0.0001 for 1000 epochs. Those hyper-parameters apply for all four datasets.

B Analysis and Ablation

This section provides in-depth analysis and explains the reasoning of the success of LaRS.

Why reasoning skill is a better guidance for demonstration selection? In TabMWP dataset, 200 examples are labeled based on the skills being showcased out of 12 manually-crafted skills labels, including "compute statistics", "compute rate of change", "Reason time schedule", "Compute probability", et. al. We investigate how the unsupervisedly discovered reasoning skills by LaRS align with human's understanding of skills. More specifically, a visualization of how human-labeled skills distribute based on the t-SNE projections of four different types of embedding is shown in Fig. 1. Both the reasoning skill encoder (reasoning skill of (Q, R)) and the reasoning policy (reasoning skill of Q) trained by LaRS demonstrate clear separation of the labeled 12 skills. At the mean time, the human-labeled skills are not well-separated by raw question embedding, and even raw rationale embeddings. This indicates that the discovered reasoning skills aligns well with human-labeled skills even without explicit labels being provided during the training. This sheds the light on why the demonstration selection based on similar reasoning skills can improve the CoT prompting.



Figure 1: t-SNE projections of reasoning skills predicted from (Q, R) (top-left), reasoning skills predicted from Q (top-right), raw question embedding (bottom-left), and raw rationale embedding (bottom-right). The 12 different colors correspond to 12 skill labels provided by human.



Random Retrieval-Q LaRS Number of in-context examples

(a) The accuracy of Random, Retrieval-Q, and, LaRS based on three different pre-trained embedding models.

(b) The accuracy of Random, Retrieval-Q, and LaRS using different number of in-context examples.

Figure 2: Performances of three different selection methods under (a) different pre-trained embedding models, and (b) different number of in-context examples.

Robustness to different pre-trained embedding models. Fig. 2a compares the performances of Random, Retrieval-Q, and LaRS based on three pre-trained embedding models, including Sentence-BERT (Reimers and Gurevych, 2019), Deberta-v2-xlarge, and, text-embedding-ada-02 (Neelakantan et al., 2022) from OpenAI. We observe that the performances of retrieval-based selection methods monotonously improve with more capable pre-trained embedding models. However, our LaRS shows consistent improvements over Retrieval-Q given the same embedding models.

Robustness to *k***: the number of in-context examples.** This study compares three selection methods, including Random, Retrieval-Q, and LaRS under three different number of in-context examples 2, 4, and 8. The results are summarized in Fig. 2b. While the accuracy monotonously improves with the increasing number of in-context examples, LaRS consistently outperforms Retrieval-Q.

C Case Study

953

954

957

960

961

962

963

964

965

967

968

970

971

972

973

974

976

977

979

To explore the examples categorized as distinct skills within the learned latent reasoning skill representation, we employed K-means clustering on the latent reasoning skills of 1,000 examples from the **TabMWP** dataset. The centroids of these clusters are detailed in Table 3. The analysis presented in this table reveals that our method effectively discerns examples showcasing specific skills, such as "Searching minimum/maximum" and "Computing rate change".

D Theoretical Analysis

To prove Theorem 1, we start with the equation of rationale generation via CoT prompting, employing the skill-based demonstration selection method denoted as g_{skill} . The process can be formalized as follows:

$$P_M(R \mid Q, g_{skill}) = \int_{\mathcal{X}^k} P_M(R \mid pt) \Pi_{i=1}^k [g_{skill}(Q_i, R_i \mid Q) d(Q_i, R_i)]$$
(5)

where Equation 5 is integrated by substituting $pt = (Q_1, R_1, \dots, Q_k, R_k, Q)$ as outlined in Equation 3, leading to:

$$P_M(R \mid Q, g_{skill}) = \int_{\mathcal{Z}} P_M(R \mid z, Q) P_M(z \mid Q) \prod_{i=1}^k [P_{skill}(z \mid Q)] dz$$
(6)

In this context, $P_{skill}(z \mid Q)$ is defined as:

$$P_{skill}(z \mid Q) = \int_{(Q',R')\in\mathcal{X}} P_M(z \mid Q',R')g_{skill}(Q',R' \mid Q)d(Q',R')dz'$$
(7)

Substituting the Definition 1 into Equation 7, leading to:

$$P_{skill}(z \mid Q) = \int_{(Q',R')\in\mathcal{X}} \int_{z'\in\mathcal{Z}} P_M(z \mid Q',R') P_E(Q',R' \mid z') P_E(z' \mid Q) dz'$$
(8)

Applying Assumption 2 into the above equation, replacing $P_M(z \mid Q', R')$ with $P_E(z \mid Q', R')$:

$$P_{skill}(z \mid Q) = \int_{(Q',R')\in\mathcal{X}} \int_{z'\in\mathcal{Z}} P_E(z \mid Q', R') P_E(Q', R' \mid z') P_E(z' \mid Q) dz'$$
$$= \int_{z'\in\mathcal{Z}} \delta(z = z') P_E(z' \mid Q) dz'$$
$$= P_E(z \mid Q)$$
(9)

By reintegrating the derived expression for $P_{skill}(z \mid Q)$ back into Equation 6, we arrive at:

$$P_M(R \mid Q, g_{skill}) = \int_{\mathcal{Z}} P_M(R \mid z, Q) P_M(z \mid Q) \prod_{i=1}^k [P_E(z \mid Q)] dz$$
(10)

Take the limit of $k \to \infty$, above equation siplifies to:

$$P_M(R \mid Q, g_{skill}) = \int_{\mathcal{Z}} P_M(R \mid z, Q) P_E(z \mid Q) dz$$
(11)

Applying Assumption 2 into the above equation, replacing $P_M(R \mid z, Q)$ with $P_E(R \mid z, Q)$:

$$P_M(R \mid Q, g_{skill}) = \int_{\mathcal{Z}} P_E(R \mid z, Q) P_E(z \mid Q) dz = P_E(R \mid Q)$$
(12)

According to Assumption 1, the example bank can approximate expert rationale generation, or $P_E(R \mid Q) = P^*(R \mid Q)$, we then conclude:

$$P_M(R \mid Q, g_{skill}) = P^*(R \mid Q) \tag{13}$$

Equation 13 means that the CoT prompting under the skill-based demonstration selection method give the
 optimal conditional distribution of rationales given questions by Definition 2. This proves the Theorem 1
 under Assumption 1 and Assumption 2.

Cluster ID	Table	Question	Skill
0	[TITLE]: School play committees Committee Boys Girls Casting 17 5 Set design 14 17 Lighting 20 20 Costume 7 4 Music 2 13	Some students at Dayton Middle School signed up to help out with the school play. Which committee has the most boys? Options: (A) set design (B) lighting (C) casting (D) costume	Search minimum/maximum
1	[TITLE]: Pairs of shoes per store Stem Leaf 1 9 2 3, 3 3 0, 2 4 2, 4 5 5, 7 6 2, 5 7 7 8 0, 2, 4, 4 9 0, 0	Ivan counted the number of pairs of shoes for sale at each of the shoe stores in the mall. How many stores have exactly 23 pairs of shoes?	Search tree leaves
2	[TITLE]: None piece of licorice \$0.07 gum drop \$0.05 gumball \$0.08 cinnamon candy \$0.01 peppermint candy \$0.08 lemon drop \$0.07	Derek has \$0.06. Does he have enough to buy a piece of licorice and a cinnamon candy? Options: (A) yes (B) no	Compute money cost
3	[TITLE]: None Number of offices Number of chairs 1 2 2 4 3 6 4 8 5 ?	Each office has 2 chairs. How many chairs are in 5 offices?	Multiplication
4	[TITLE]: None popcorn balls \$1/kilogram coffee cake \$3/kilogram blueberry bars \$2/kilogram cream cheese bars \$2/kilogram lemon bars \$3/kilogram	Sarah went to the store and bought 2 kilograms of blueberry bars. How much did she spend? (Unit: \$)	Compute money cost
5	[TITLE]: None x y 12 19 13 9 14 2	The table shows a function. Is the function linear or nonlinear? Options: (A) linear (B) nonlinear	Compute rate of change
6	[TITLE]: Tractors Farmer Number of tractors Farmer Judy 4 Farmer Joe 7 Farmer Megan 7 Farmer Rick 4 Farmer Jane 4	Some farmers compared how many tractors they own. What is the mode of the numbers?	Compute statistics
7	[TITLE]: None pink sweater \$6.69 pair of brown pants \$9.66 plaid scarf \$2.45 pair of sandals \$7.69 white polo shirt \$4.86	How much money does Heather need to buy a pair of brown pants and a plaid scarf? (Unit: \$)	Compute money cost
8	[TITLE]: Tour bus schedule Location Arrive Depart the riverfront 9:55 A.M. 10:20 A.M. the zoo 10:35 A.M. 11:30 A.M. art museum 12:05 P.M. 12:30 P.M. science museum 1:00 P.M. 1:45 P.M. skyscraper 1:50 P.M. 2:20 P.M. governor's mansion 2:50 P.M. 3:45 P.M. old building 4:00 P.M. 4:45 P.M. famous bridge 5:15 P.M. 5:40 P.M. the aquarium 6:20 P.M. 7:00 P.M. landmark sculpture 7:45 P.M. 8:20 P.M.	Look at the following schedule. Which stop does the bus depart from at 11.30 A.M.? Options: (A) zoo (B) riverfront (C) old building (D) science mu- seum	Reason time schedule

Cluster ID	Table	Question	Skill
9	[TITLE]: None poppyseed muffin \$2.31 bowl of yogurt \$1.35 blueberry pancakes \$7.28 hash browns \$4.56 bowl of granola \$2.97 bagel with cream cheese \$2.56	Max has \$13.33. How much money will Max have left if he buys a bagel with cream cheese and blueberry pancakes? (Unit: \$)	Compute money cost
10	[TITLE]: Balloons sold Day Number of balloons Wednesday 568 Thursday 586 Friday 558 Saturday 565	The manager of a party supply store researched how many balloons it sold in the past 4 days. On which day did the store sell the most balloons? Options: (A) Wednesday (B) Thursday (C) Friday (D) Saturday	Search minimum/maximum
11	[TTTLE]: None forklift \$9,987.00 dump truck \$9,543.00 race car \$8,370.00 crane \$6,996.00 bulldozer \$7,547.00 hydrofoil \$8,047.00	How much more does a forklift cost than a dump truck? (Unit: \$)	Compute money cost

Table 3: The closest examples to the 12 cluster centers computed by K-Means clustering method on reasoning skill latent variables.